

# GPU Basics

## Category: Pleiades

A graphics processing unit (GPU) is a hardware device that may be able to accelerate an algorithm or computer code. If you want general information on GPUs, you might start with the Wikipedia article "[GPGPU](#)" and the [GPGPU website](#), which has information for developers.

If you have an application that can take advantage of GPU technology, you can use the GPU nodes on Pleiades. Specifically, the 64 Pleiades Westmere nodes on rack 219 (r219i[0-3]n[0-15]) include one NVIDIA Tesla M2090 GPU per node. Each M2090 computing module comprises a computing subsystem with a Tesla 20-series GPU and high speed memory, and is connected to the Westmere node via a PCI Express bus.

The [NVIDIA Developer Zone](#) has specific information about NVIDIA GPUs and the programming model for them.

To use the Westmere+GPU nodes, specify the processor model type `model=wes_gpu` in your PBS script:

```
#PBS -l select=xx:ncpus=yy:model=wes_gpu
```

Submit your PBS jobs to the **gpu** queue as follows:

```
% qsub -q gpu job_script
```

To check the status of your jobs submitted to the **gpu** queue:

```
% qstat gpu -u your_username
```

Get basic hardware information about the Pleiades GPUs, as follows:

```
pfe20% qsub -I -q gpu
r219i0n0% /usr/bin/nvidia-smi -q
and/or
r219i0n0% module load comp-pgi/11.6
r219i0n0% pgaccelinfo
```

The output from **pgaccelinfo** shows:

```
CUDA Driver Version:           4000
NVRM version: NVIDIA UNIX x86_64 Kernel Module  275.09.07  Wed Jun  8 14:16:46 PDT 2011

Device Number:                  0
Device Name:                    Tesla M2090
Device Revision Number:         2.0
Global Memory Size:             5636554752
Number of Multiprocessors:      16
```

Number of Cores:	512
Concurrent Copy and Execution:	Yes
Total Constant Memory:	65536
Total Shared Memory per Block:	49152
Registers per Block:	32768
Warp Size:	32
Maximum Threads per Block:	1024
Maximum Block Dimensions:	1024, 1024, 64
Maximum Grid Dimensions:	65535 x 65535 x 65535
Maximum Memory Pitch:	2147483647B
Texture Alignment:	512B
Clock Rate:	1301 MHz
Initialization time:	20609 microseconds
Current free memory:	5552726016
Upload time (4MB):	1547 microseconds ( 761 ms pinned)
Download time:	1113 microseconds ( 681 ms pinned)
Upload bandwidth:	2711 MB/sec (5511 MB/sec pinned)
Download bandwidth:	3768 MB/sec (6159 MB/sec pinned)

To use the Pleiades GPUs you have two possibilities:

1. The PGI accelerator model requires you to annotate your source code (Fortran or C) with directives describing sections of your code that are to be executed on the GPU. You need to use one of the PGI compiler modules (for example, **module load comp-pgi/11.6**). Refer to the [Portland Group website](#) for information on their accelerator compilers.
2. You may write (or rewrite) portions of your code in CUDA (Compute Unified Device Architecture; see the [NVIDIA Developer Zone](#)) or in CUDA Fortran. If you choose the latter, you will need to use a PGI compiler module (such as **module load comp-pgi/11.6**). If you select the former, you will need to load a **cuda** module (such as **module load cuda/4.0**) and use the CUDA tools. You will also need to determine how to link your new CUDA code with the rest of your program. If you need assistance, contact the [NAS Control Room](#), and our consultants will lend a hand.

Currently, no direct communication exists between a GPU on one node and a GPU on another node. If such communication is required, the data must go via the PCI Express bus from the GPU to the Westmere CPU and via MPI from one CPU to another.

If you delve into CUDA programming, the following book may be useful:

[CUDA by Example: An Introduction to General-Purpose GPU Programming](#), Jason Sanders and Edward Kandrot

Computing at NAS -> Computing Hardware -> Pleiades -> GPU Basics  
<http://www.nas.nasa.gov/hecc/support/kb/entry/298/?ajax=1>